



# Neue Vorschläge zur Evidenzbewertung

Stefanie Thomas & Ralf Bender



# Neue Vorschläge zur Evidenzbewertung ...

## Methodendiskussionen im IQWiG



- SGB V 139a: „Bewertung des medizinischen Nutzens nach den international **anerkannten Standards** der evidenzbasierten Medizin“
- 2008: Start **IQWiG im Dialog** über **methodische Herausforderungen**
- Start mit „Replikation“; weiter: Bewertung von **Schaden**; **Klinische Relevanz**; **Heterogenität / Subgruppen**; **Unsicherheit (Anhaltspunkt!)**; **Bedeutung der Zulassung**; **Studien mit erlaubtem Behandlungswechsel**
- 2015: **Wie konfirmatorisch ist HTA?**  
→ HTA ist ein geordneter Prozess, ein konfirmatorischer Nachweis (der **Wirksamkeit**) die Eintrittspforte für ... Nutzenbewertung.

## 2015: Kritik aus der „Real World“..

- **2015, Herbstsymposium: Real World Data**
- „Not all relevant questions can be addressed by real-world data today but [...] **non-RCT information will become increasingly relevant for assessment of benefits, risks, comparative effectiveness, and value.**“ (Eichler 2015)
- „Für die GKV sind randomisierte klinische Studien der **Goldstandard, aber die Wissenschaft denkt schon weiter.**“ (DAZ 2016)



## 2015: Kritik aus der „Real World“..

- 2015, **Herbstsymposium: Real World Data**
- „Not all relevant questions can be addressed by real-world data today but [...] **non-RCT information will become increasingly relevant for assessment of benefits, risks, comparative effectiveness, and value.**“ (Eichler 2015)
- „Für die GKV sind randomisierte klinische Studien der Goldstandard.“ (DAZ 2016)



### Fazit

- Einigkeit – **RCT am besten geeignet für valide Nutzensaussagen**
  - **Non-RCT vielleicht ergänzend geeignet** (wie genau, blieb offen)
- „Erreichte Standards werden aufgegeben, und **ohne Not soll eine erhöhte Unsicherheit in Kauf** genommen werden. [...].  
Hier müssen wir gegensteuern.“ (Windeler 2015)

## Heutiges Thema: Kritik an zu unkritischen „Standards“



- **“The knowledge system underpinning healthcare is not fit for purpose and must change. The medical literature is biased and inundated with poor quality trials.”**  
(Roberts 2015)
- „The main utility of systematic reviews has been to reveal how miserably unreliable biomedical evidence is. [...] Given the prestige that systematic reviews have acquired in the hierarchy of evidence, more and more people do more of them, creating bubbles of reviews on largely unreliable evidence with more than deserved credibility assigned to them.” (Ioannidis 2015)
- **Kritik an mangelhaften Studien, Publikationskultur und (dafür) unangemessenen Review-Methoden → Aufruf zur Änderung!**

# Überblick: Vorschläge zur strengeren Evidenzbewertung

## (1) Umgang mit Publication Bias

- Einschluss nur prospektiv registrierter Studien

## (2) „Risk of fraud ... doubtful trials“: Validitätsbewertung

- Härtere Überprüfungen von Design und Daten
- Malus bei Interessenskonflikten

## (3) „Revised standards for statistical evidence“

- Fallzahlplanung („DARIS“) für Reviews
- p-Werte adjustiert für Multiplizität (Sequenz; Endpunkte)
- Bayes-Faktoren
- strengeres Signifikanzniveau (0,005 oder sogar 0,001)

## (4) „Assessment of clinical significance“

- Abgleich mit (Ir)Relevanzschwellen der Fallzahlplanung

# (1) „Filter Failure“

(Ben Goldacre 2015)

## ■ Hintergrund: Publikation und Reporting-Bias

“...overestimation of efficacy and the underestimation of safety risks of interventions. ... Mandatory prospective registration of trials and public access to study data via results databases need to be introduced on a worldwide scale.” (z. B. McGauran et al. 2010)

## ■ Bekannte Gegenmaßnahmen

- ICMJE-Statement 2004 (De Angelis 2004)
- Verpflichtende Registrierung seitens FDA/EMA
- **...reichen nicht:** (Prayle 2012, Anderson 2015, Goldachre 2015)
  - mangelnde Stringenz, Audits bzw. Sanktionierung
  - Zusammenhang mit (AM-) Zulassung
  - nicht retrospektiv wirksam
  - Einsichtnahme...
- **„Despite decades of exhortation about trial publication, about half of all trials are unpublished.“** (Roberts 2015)

# Bessere „Filter“?

## Gegenmaßnahmen auf Seiten der Reviewer:

- **Ausschluss nicht registrierter Studien** (ab 2010 publizierte)
- „Including only prospectively registered trials in systematic reviews will improve validity and readability.“ (Roberts 2015)
- Diskussion im Gange:
  - „If the best evidence is not prospectively registered the validity will be reduced by excluding it.“ (Alper 2015)
  - „[...] we should not confuse issues of identifying all available data [...] and assessing the validity of that data.“ (Tovey 2015)



# Informationsbeschaffung ↔ Bewertung fehlender Information

- **IQWiG – Allgemeine Methoden**
  - Probleme sowohl bezüglich Publikations-, als auch Reporting Bias!
  - reguläre Recherchen in Studienregistern
  - Anfragen bei Herstellern (und Autoren)
  - „bekanntermaßen“ unvollständige Daten: → **Fazit unter Vorbehalt**
- **Probleme insbesondere bei Methoden / Medizinprodukten (NMV)**
  - N05-03C Stammzelltransplantation (2011)
    - nur Daten von 2 von 5 Studien zu relevanten Vergleichen zur Verfügung – Anfragen erfolglos
  - D12-01 Kontinuierliche interstitielle Glukosemessung (2015)
    - (S)UE grob unvollständig berichtet

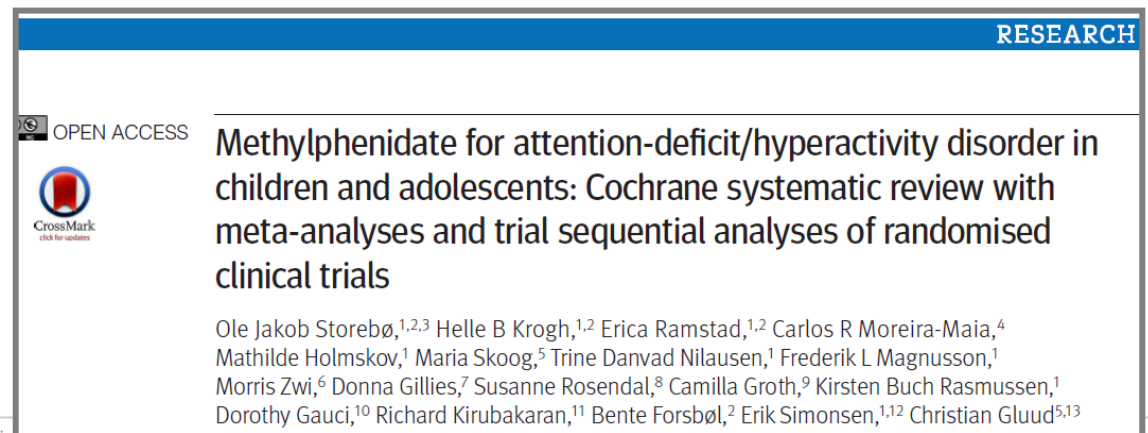
## (2) Validitätsbewertung bei Verdacht..

1. „The risk of fraud cannot be ignored.“ (Roberts 2015)

- **Statistische Checks durchführen**
- **Originaldaten anfordern**
- **Studien bei fehlender Antwort ausschließen**


2. **Sponsoring und Interessenskonflikte**

- „The risk of bias domains [...] and industry funding have been shown to be of **particular importance**.“ (Jakobsen 2015; Lundh 2012)
- Neue Standard-Domäne zur Bewertung des Risk of bias:  
„**Vested interest**“ (Storebø 2015)



RESEARCH

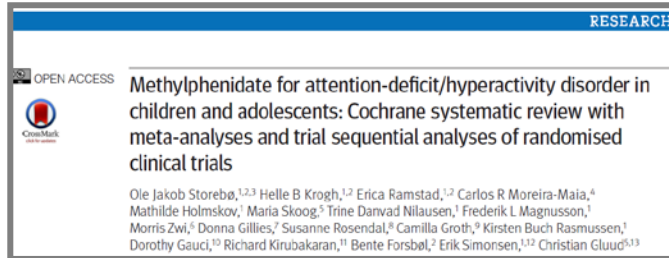
OPEN ACCESS

 CrossMark  
click for updates

**Methylphenidate for attention-deficit/hyperactivity disorder in children and adolescents: Cochrane systematic review with meta-analyses and trial sequential analyses of randomised clinical trials**

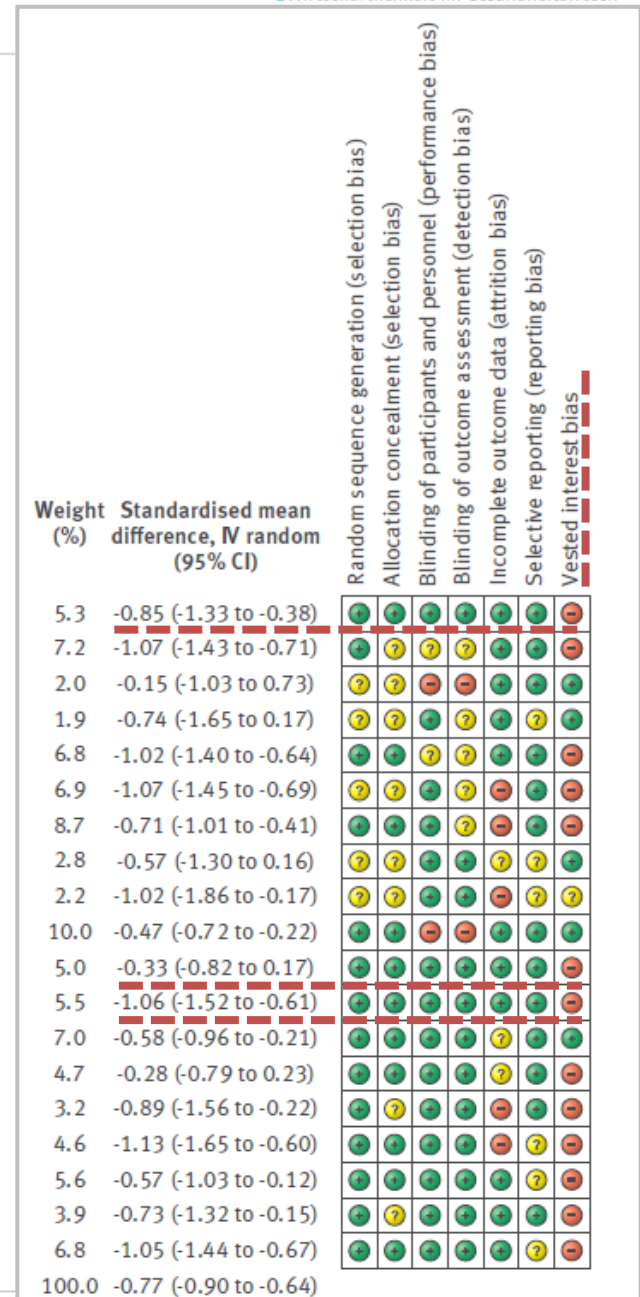
Ole Jakob Storebø,<sup>1,2,3</sup> Helle B Krogh,<sup>1,2</sup> Erica Ramstad,<sup>1,2</sup> Carlos R Moreira-Maia,<sup>4</sup> Mathilde Holmskov,<sup>1</sup> Maria Skoog,<sup>5</sup> Trine Danvad Nilausen,<sup>1</sup> Frederik L Magnusson,<sup>1</sup> Morris Zwi,<sup>6</sup> Donna Gillies,<sup>7</sup> Susanne Rosendal,<sup>8</sup> Camilla Groth,<sup>9</sup> Kirsten Buch Rasmussen,<sup>1</sup> Dorothy Gauci,<sup>10</sup> Richard Kirubakaran,<sup>11</sup> Bente Forsbøl,<sup>2</sup> Erik Simonsen,<sup>1,12</sup> Christian Gluud<sup>5,13</sup>

# Ein Beispiel: Storebø et al. 2015



## Erweitertes Risk of Bias – Tool

- „vested interest“
- inadäquat → „high risk of bias“
- **Konsequenzen für IQWiG-Berichte?**
- **Aktueller Standard: Bewertung allein von Designaspekten und Ergebnisberichten**



### (3) „Revised Standards for Statistical Evidence“?

#### Kritik an Methodik für systematischen Reviews:

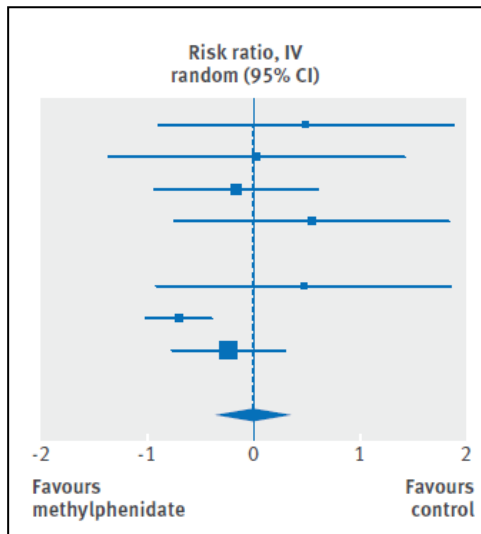
- „nur retrospektive Beobachtungsstudie... “ (*...wie konfirmatorisch?*)
- **konsekutives Studien-Hinzufügen** in Meta-Analysen  
↔ sequenzielles Testen
- **2 Vorschläge** (Jakobson 2014):
  - Sequenzielle Methoden für Meta-Analysen  
→ **strengere Tests, statistisch valide(re) Ergebnisse**
  - Korrespondierende Fallzahlberechnung  
→ **„bessere Interpretation nicht-signifikanter Ergebnisse“**
- ⇒ **„Trial Sequential Analysis (TSA) und Diversity adjusted required informations size (DARIS) (Wetterslev 2008)**
- **...1 Beispiel** zur Visualisierung der „DARIS“ (...Storebø 2015)

# DARIS am Fallbeispiel

## Storebø 2015

- Endpunkt: SUE
- n=185 Studien
- Berichtet nur in 9 (5%) Studien (1.721 Patienten)
- Studien „high risk of bias“.

...konventionell:



Neu:  
**DARIS =**  
21.593 Patienten  
„type II error cannot be excluded.“  
(Keine TSA-Anwendung für sign. Ergebnisse)



We conducted trial sequential analysis on the “total serious adverse events” outcome, involving nine parallel group trials. We had planned to use a relative risk reduction of 20% but owing to too large a distance between the accrued information and the required information the program rejected to calculate and draw an interpretable figure. We therefore increased the relative risk reduction to 25%. We included trials with zero serious adverse events by substituting a constant of 0.5 for zero. We calculated the diversity

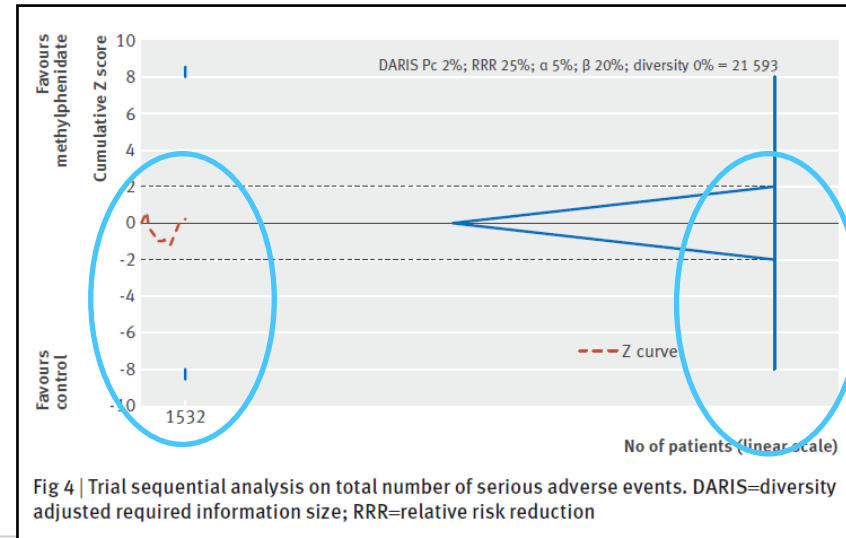


Fig 4 | Trial sequential analysis on total number of serious adverse events. DARIS=diversity adjusted required information size; RRR=relative risk reduction

# Trial Sequential Analysis – DARIS

- **Laufende Diskussionen:**
  - **Methodisch** noch nicht ausgereift?  
„...we suggest that a routine requirement for such a calculation might cause more difficulties than it solves.“ (Kulinskaya 2015)
  - Ziel von **DARIS**: Nutzen oder Schaden als **Interpretationshilfe**?
  - ...welches „**Fallzahldelta**“ ...?
  - Ziel von **TSA**: Notwendigkeit der **Multiplizitätskontrolle**?
- **Aktueller IQWiG-Standard: Lokale Konfidenzintervall-Schätzung.**
  - **pro Endpunkt...**

## Multiplizität, contd. ...

### Weiterer Kritikpunkt **Multiplizität durch $n > 1$ Endpunkte (EP)**

- „Pragmatischer“ Lösungsvorschlag:  
„...**dividing the pre-specified P-value threshold** with the value **halfway between 1 (no adjustment) and the number of primary outcome comparisons (Bonferroni adjustment).**“ (Jakobson 2015)
- Also: Divisor bei 2 EP: 1,5 (statt 2); bei 3 EP: 2 (statt 3) usw. ...

### ▪ **Problembewusstsein im IQWiG...**

„**Problem in systematischen Übersichten nicht komplett lösbar.**“

- **Allein 3 Endpunktkategorien** mit regelmäßig  $>1$  berichteten Endpunkten mit  $>1$  Operationalisierungen
- **Konsistenzforderung** (z. B. innerhalb EP und EP-Kategorie)
  - Keine Adjustierung für multiples Testen für die endpunktübergreifende Bewertung
  - **Forderungen von Daten pro Kategorie!**



## Statisticians issue warning on *P* values

Statement aims to halt missteps in the quest for certainty. Baker 2016

### Warnung I:

- **Angabe und / oder Interpretation des p-Wertes problematisch:**  
„Even a low P-value from a meta-analysis can be misleading. ”
- **„Kompatibilität der Alternative  $H_1$  mit den Daten prüfen!“** (Jakobson 2015)

Bayes Faktor: 
$$BF_{01} = \frac{P(\text{Daten beob.} | \text{Modell } H_0)}{P(\text{Daten beob.} | \text{Modell } H_1)}$$
 (adapt. n. Berry 1996)

→ „Verhältnis der Wahrscheinlichkeit der beobachteten Daten, wenn die Nullhypothese  $H_0$  wahr ist und der, wenn Alternative  $H_1$  wahr ist.“

- **Ziel: Angemessenere Ergebnisinterpretation durch (ergänzende) Betrachtung des Bayes Faktors.**



# Bayes-Faktoren die bessere Alternative?

- Vor der Anwendung wären da noch...

## 1. Die Alternative?

- „Bayes factor can be defined differently...“
- Spezifikation der (!) Alternative  $H_1$  dafür unabdingbar
  - „However, BF will still be misleading when an **unrealistic large anticipated intervention effect** is confirmed by ‚play of chance‘ by an unrealistically large observed effect...“

## 2. Und welche Niveaus würden für Bayes-Faktoren gelten?

P-Wert	BF <sub>01</sub>	BF <sub>10</sub>	⇒ Ideen für das „Fazit“
<0,05	>1	<1	⇒ „ <b>more credit</b> to the null hypothesis being true.“
(<0,05)	„high“	(„low“)	⇒ „ <b>interpretation of results with caution...</b> effect lower than anticipated..“
(<0,05)	<0,1	>10	⇒ „... <b>may be chosen as threshold for significance</b> .“

(nach Jakobson 2015)

## .. Oder strengere Anforderung an die p-Werte?

### Warnung II:

- „Significance tests at inappropriately high levels of significance... A root cause of nonreproducibility.“ (Johnson 2013)

### Vorschlag:

1. Korrespondierende p-Werte und Bayes-Faktoren
2. Niveaus für „signifikante“ Bayes-Faktoren an p-Werte anlegen
3. ...neue p-Werte!

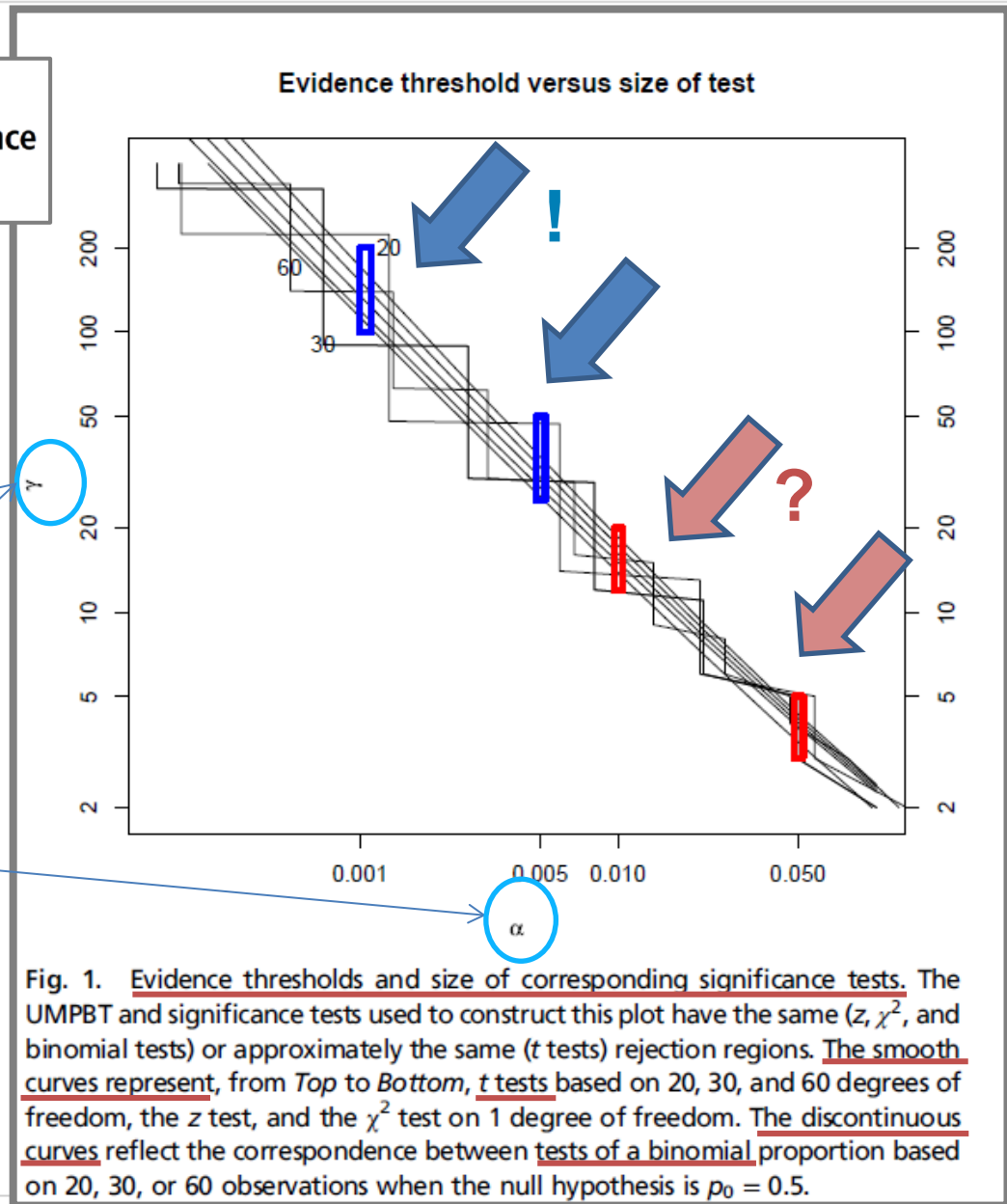


# .. Oder strengere Anforderung an die p-Werte?

**PNAS**  
**Revised standards for statistical evidence**  
Valen E. Johnson<sup>1</sup>  
Department of Statistics, Texas A&M University, College Station, TX 77843-3143

„Evidence threshold“  
(Bayes Factor)

„Size of significance test“  
(p-Werte)



## .. Oder strengere Anforderung an die p-Werte?

### Warnung II:

- „Significance tests at inappropriately high levels of significance... A root cause of nonreproducibility.“ (Johnson 2013)

### Vorschlag:

1. Korrespondierende p-Werte und Bayes-Faktoren
2. Niveaus für „signifikante“ Bayes-Faktoren an p-Werte anlegen
3. **Konsequenz: Einforderungen 10-fach strengerer p-Werte.**

BF <sub>10</sub>		Bewertung BF	P-Wert
<10 <sup>0</sup>		„negative“ (pro Nullhypothese)	
<10 <sup>0</sup> – 10 <sup>1/2</sup>	(1 – 3)	„barely worth mentioning“	
10 <sup>1/2</sup> – 10 <sup>1</sup>	(3 – 10)	„positive“ / „substantial“	
10 <sup>1</sup> – 10 <sup>3/2</sup>	<b>(10 – 32)</b>	<b>„strong“</b>	<b>0,005</b>
10 <sup>3/2</sup> – 10 <sup>2</sup>	<b>(32 – 100)</b>	<b>„very strong“</b>	
> 10 <sup>2</sup>	<b>&gt;100</b>	<b>„decisive“</b>	<b>0,001</b>

nach Jeffreys 1961

## Last but not least: „Die“ Relevanzbewertung

„More than statistical significance is required before an intervention is declared as being effective...“

"Klinische Studien werden mit dem Ziel durchgeführt, Ergebnisse zu liefern, die in der Praxis zu einer Verbesserung der Patientenversorgung führen. Allgemein wird jedoch beklagt, dass die Charakterisierung von Studienergebnissen als 'signifikant' oder 'nicht signifikant' keine Aussagen über ihre klinische Relevanz erlaubt."

Windeler J, Conradt C. Wie können "Signifikanz" und "Relevanz" verbunden werden?. Med Klin 1999; 94: 652-5.

**Stefan Lange, IQWiG im Dialog 2010**

## Last but not least: „Die“ Relevanzbewertung

„More than statistical significance is required before an intervention is declared as being effective...“

- When **surrogate outcomes** or **continuous outcomes** are used to assess intervention effects, it is often unclear if a given statistical significant effect has any patient relevant clinical significance.
- „Relating trial results to the ‘**minimal relevant clinical ...difference**’ used to calculate the predefined sample size as well as calculating Bayes factor based on this ‘minimal relevant clinical difference’, provide indications about the clinical significance of intervention effects.“
- „To assess the clinical significance of intervention effects it is important to perform a **thorough assessment of the balance** between beneficial and harmful effects.“

(Jakobson 2015)

- „Ein weites Feld...“

(Vater Briest, 1896)

## Last but not least: „Die“ Relevanzbewertung

- **Differenzierung der Ebenen im IQWiG:**
  - **Individuelle Veränderungen  $\neq$  Gruppenunterschiede  $\neq$  endpunktübergreifende Abwägung**
  - **Irrelevanzgrenzen  $\neq$  Relevanzgrenzen**
  
- ⇒ Nachweis definierter **Effektstärken bei PRO / Skalen** und **Responderanalysen mittels individueller MID** [„Minimal Important Difference“]
- ⇒ Gestaffelter Effektstärkennachweis : **Ausmaß-Methodik**
  
- **Patientenrelevante Endpunkte und valide Surrogate**
  
- **Überblick zur Abwägung: Landkarte der Beleglage seit 2010**  
...und zu Ausmaß und Wahrscheinlichkeit des Zusatznutzens seit 2011

# Last but not least: „Die“ Relevanzbewertung

## ..wenige internationaler Standards:

- Indikationsbezogene Endpunktdiskussionen...
  - **PROs nur in verblindeten Studien** aussagekräftig! (FDA 2009)
- **Für die Abwägung:** Benefit-Risk Methodology Project (EMEA 2009-EMA 2014)
  - „...recommendation [...] the use of **Effects Table (ET)** as a tool to summarise the key benefits and risks.“



# Zu streng, zu lasch.. Was ist angemessen?

## (1) Reaktion auf Publication- & Reporting Bias

- Einschluss nur prospektiv registrierter Studien?

## (2) Strengere Studienbewertung

- Härtere Überprüfungen von Daten?
- Malus bei Interessenskonflikten?

## (3) Strengere statistische Standards

- Fallzahlplanung für Reviews?
- p-Werte adjustiert für multiple Multiplizität?
- höheres Signifikanzniveau (0,005 oder sogar 0,001)?
- Bayes-Faktoren als Interpretationshilfen?

## (4) Bewertung der klinischen Relevanz

- Abgleich mit (Ir)Relevanzschwellen der Fallzahlplanung?
- Verfahren für eine Abwägung?



## Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen (IQWiG)

- Im Mediapark 8
- D-50670 Köln
  
- Telefon +49-221/3 56 85-0
- Telefax +49-221/3 56 85-1
  
- [info@iqwig.de](mailto:info@iqwig.de)
- [www.iqwig.de](http://www.iqwig.de)



„Eine Hauptursache der Armut in den Wissenschaften ist meist eingebildeter Reichtum. Es ist nicht ihr Ziel, der unendlichen Weisheit eine Tür zu öffnen, sondern eine Grenze zu setzen dem unendlichen Irrtum.“  
(Bertolt Brecht 1943)

# Literatur

- Alper BS. Re: The knowledge system underpinning healthcare is not fit for purpose and must change. BMJ online 2015. URL: <http://www.bmj.com/content/350/bmj.h2463/rapid-responses>
- Anderson ML, Chiswell K, Peterson ED, Tasneem A, Topping J, Califf R. Compliance with Result Reporting at ClinicalTrials.gov. NEJM 2015; 372(11) 1031 - 1039
- Baker M. Statisticians issue warning on P values. Nature 2016; 531 151
- Berry DA. Statistics A Bayesian Perspective. Duxbury Press 1996
- Brecht B. Leben des Galilei. Frankfurt: Suhrkamp. Uraufführung, 1. Version 1943, Schauspielhaus Zürich
- DAZ.online. 19. Eppendorfer Dialog, RCT - Goldstandard oder Auslaufmodell? URL: <https://www.deutsche-apotheker-zeitung.de/news/artikel/2016/04/21/rct-goldstandard-oder-auslaufmodell>
- De Angelis C, Drazen JM, Frizelle FA, Haug C, Hoey J, Horton R, Kotzin S, Laine C, Marusic A, Overbeke AJ, Schroeder TV, Sox HC, Van Der Weyden MB. Clinical trial registration: a statement from the International Committee of Medical Journal Editors. Ann Intern Med. 2004;141(b): 477–78..
- Eichler H-G, Baird LG, Barker L et al. From Adaptive Licensing to Adaptive Pathways: Delivering a Flexible Life-Span Approach to Bring New Drugs to Patients. Clin Pharmacol Ther 2015; 79(3) 234-246
- EMA 2014. Benefit-risk methodology project Update on work package 5: Effects Table pilot. URL: [http://www.ema.europa.eu/docs/en\\_GB/document\\_library/Report/2010/04/WC500089603.pdf](http://www.ema.europa.eu/docs/en_GB/document_library/Report/2010/04/WC500089603.pdf)
- FDA 2009. Guidance for industry: patient-reported outcome measures; use in medical product development to support labeling claims. URL: <http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/UCM193282.pdf>
- Fontane T. Effi Briest. Reclam Verlag 1986
- Goldacre B. Perfecting Evidence-based Medicine. Cochrane Kolloquium 2015. URL: [https://www.youtube.com/watch?v=h2giHu8pHFY&list=PLCo8P5\\_ppmQgoKI5ofhvBn-0yZnylWoMD&index=2](https://www.youtube.com/watch?v=h2giHu8pHFY&list=PLCo8P5_ppmQgoKI5ofhvBn-0yZnylWoMD&index=2)
- Ioannidis JPA. Can too many systematic reviews and meta-analysis do harm? Cochrane Colloquium 2015. URL: [https://www.youtube.com/watch?v=h2giHu8pHFY&list=PLCo8P5\\_ppmQgoKI5ofhvBn-0yZnylWoMD&index=2](https://www.youtube.com/watch?v=h2giHu8pHFY&list=PLCo8P5_ppmQgoKI5ofhvBn-0yZnylWoMD&index=2)
- IQWiG Allgemeine Methoden Vs. 4.2. 2015. URL: [https://www.iqwig.de/download/IQWiG\\_Methoden\\_Version\\_4-2.pdf](https://www.iqwig.de/download/IQWiG_Methoden_Version_4-2.pdf)
- Jakobsen, J.C., Wetterslev, J., Winkel, P., Lange, T. & Gluud, C.: Thresholds for statistical and clinical significance in systematic reviews with meta-analytic methods. BMC Med. Res. Methodol. 2014; 14, 120.
- Jeffreys H (1961) Theory of Probability (Oxford Univ Press, Oxford), 3rd Ed.
- Johnson, V.E.: Revised standards for statistical evidence. Proc. Natl. Acad. Sci. U.S.A. 2013; 110, 19313-19317.
- Kulinskaya Re: The knowledge system underpinning healthcare is not fit for purpose and must change. BMJ online 2015. URL: <http://www.bmj.com/content/350/bmj.h2463/rapid-responses>
- Lundh A, Sismondo S, Lexchin J, Busuioic OA, Bero L. Industry sponsorship and research outcome. Cochrane Database Syst Rev 2012
- McGauran N, Wieseler B, Kreis J, Schüler Y-B, Kölsch H, Kaiser T. Reporting bias in medical research - a narrative review. Trials 2010, 11:37
- Prayle AP, Hurley MN, Smyth AR. Compliance with mandatory reporting of clinical trial results on ClinicalTrials.gov: cross sectional study. BMJ 2011;344:d7373
- Roberts, I., Ker, K., Edwards, P., Beecher, D., Manno, D. & Sydenham, E.: The knowledge system underpinning healthcare is not fit for purpose and must change. BMJ 2015; 350, h2463.
- Storebø OJ, Krogh HB, Ramstad E, Moreira-Maia CR, Holmskov M, Skoog M, Nilausen TD, Magnusson FR, Zwi M, Gillies D, Rosendal S, Groth C, Rasmussen K, Gauci D, Kirubakaran R, Forsbøl B, Simonsen E, Gluud C. Methylphenidate for attention-deficit/hyperactivity disorder in children and adolescents: Cochrane systematic review with meta-analyses and trial sequential analyses of randomised clinical trials. BMJ 2015;351:h5203
- Tovey DI. Re: The knowledge system underpinning healthcare is not fit for purpose and must change. BMJ online 2015. URL: <http://www.bmj.com/content/350/bmj.h2463/rapid-responses> 20.
- Unternehmen heute (online) 22.04.2016 URL: <http://unternehmen-heute.de/news.php?newsid=352025>
- Wetterslev J, Thorlund K, Brok J, Gluud C. Trial sequential analysis may establish when firm evidence is reached in cumulative meta-analysis. J Clin Epidemiol 2008; 61, 64-75
- Windeler J, Herbstsymposium 2015. URL zur Pressemitteilung: <https://www.iqwig.de/de/presse/pressemitteilungen/pressemitteilungen/herbst-symposium-2015-vortrage-jetzt-online.7104.html> bzw. Einführung und Schlusswort: [https://www.iqwig.de/download/HS15\\_Einfuehrung\\_und\\_Schlusswort\\_Juergen\\_Windeler.pdf](https://www.iqwig.de/download/HS15_Einfuehrung_und_Schlusswort_Juergen_Windeler.pdf)